

AUTOMATIC EXTRACTION DEVICE FOR RELATIVE KEYWORD, DOCUMENT RETRIEVING DEVICE AND DOCUMENT RETRIEVING SYSTEM USING THESE DEVICES

Publication number: JP11025108

Publication date: 1999-01-29

Inventor: SATO MITSUHIRO; NOGUCHI NAOHIKO; SUGANO YUJI; NOMOTO MASAKO; INABA MITSUAKI; FUKUSHIGE TAKAO

Applicant: MATSUSHITA ELECTRIC IND CO LTD

Classification:

- international: G06F17/30; G06F17/30; (IPC1-7): G06F17/30

- European: G06F17/30T1E

Application number: JP19970176822 19970702

Priority number(s): JP19970176822 19970702

Also published as:



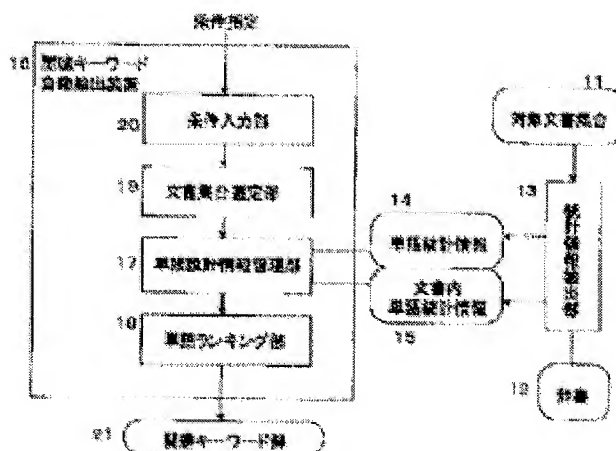
EP0889419 (A2)
US6212517 (B1)
EP0889419 (A3)
CN1206158 (A)
EP0889419 (B1)

more >>

Report a data error here

Abstract of JP11025108

PROBLEM TO BE SOLVED: To automatically extract a relative keyword which is matched with the characteristics of a document to be practically retrieved and which is capable of obtaining one or more retrieval results at the time of executing retrieval using the keyword. **SOLUTION:** An automatic extraction device for relative keywords is provided with a document set selection part 19 for specifying a partial set of each document based on the attribute information, input retrieval expression, etc., of the document, a word statistic information management part 17 for managing the statistic information of respective words in the whole objective document 11 and words appearing in each document as well as their statistic information 15; and a word ranking part 18 for calculating the importance of each word appearing in a partial set of a certain document and for aligning respective words in the order of importance, wherein the management part 17 quickly finds out the statistic information of respective words in the whole document and a specified partial set of the document. Consequently, words appearing in a certain document set can be ranked based on their importance and a part of the ranked words can be presented as a relative keyword.



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-25108

(43) 公開日 平成11年(1999) 1月29日

(51) Int.Cl.⁵
G 0 6 F 17/30

識別記号

F I
G 0 6 F 15/401 3 1 0 A
15/40 3 7 0 A

審査請求 未請求 請求項の数17 O L (全 17 頁)

(21) 出願番号 特願平9-176822

(22) 出願日 平成9年(1997) 7月2日

(71) 出願人 000005821
松下電器産業株式会社
大阪府門真市大字門真1006番地
(72) 発明者 佐 藤 光 弘
大阪府門真市大字門真1006番地 松下電器
産業株式会社内
(72) 発明者 野 口 直 彦
大阪府門真市大字門真1006番地 松下電器
産業株式会社内
(72) 発明者 菅 野 祐 司
大阪府門真市大字門真1006番地 松下電器
産業株式会社内
(74) 代理人 弁理士 蔵合 正博

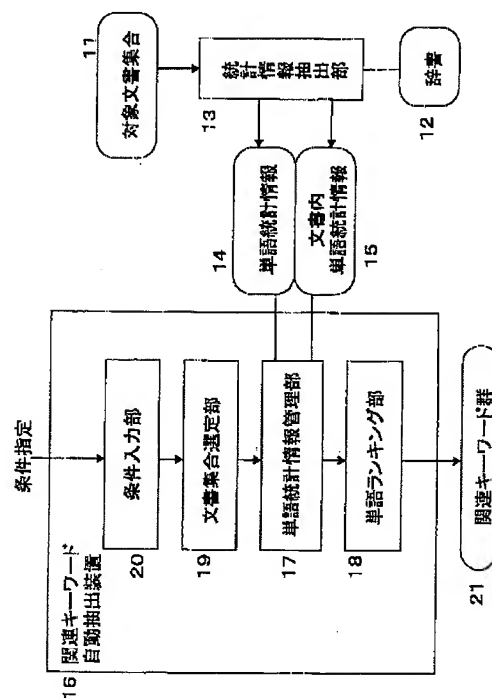
最終頁に続く

(54) 【発明の名称】 関連キーワード自動抽出装置、文書検索装置及びこれらを用いた文書検索システム

(57) 【要約】

【課題】 実際の検索対象文書の特性に即し、かつそのキーワードによる検索を実行した場合少なくとも1件以上の検索結果が得られるような関連キーワードを自動抽出すること。

【解決手段】 関連キーワード自動抽出装置として、各文書の属性情報や入力検索式などに基づいて文書の部分集合を特定する文書集合選定部19と、各単語の対象文書11全体における統計情報14および文書毎に出現する単語とその統計情報15を管理する単語統計情報管理部17と、単語統計情報14、15を基に、或る文書の部分集合に出現する各単語の重要度を算出して重要度の順に整列する単語ランキング部18とを設け、単語統計情報管理部により、文書全体、および特定された文書部分集合における各単語の統計情報を高速に求める。これにより、或る文書集合に出現する単語を、その重要度に基づいてランキングし、その一部を関連キーワードとして提示することができる。



【特許請求の範囲】

【請求項1】 辞書を用いて対象文書集合の各文書に出現する単語または単語の組の出現頻度や分布などの統計情報があらかじめ抽出されている文書集合に対して、各文書に付与された属性情報やユーザが入力した検索式などに基づいて文書の部分集合を特定する文書集合選定部と、各単語の対象文書全体における統計情報、および各文書ごとの当該文書に出現する単語とその統計情報を管理する単語統計情報管理部と、各単語の全文書および各文書ごとの統計情報を基に、特定された部分集合に出現する各単語の重要度を算出して重要度の順に整列する単語ランキング部とを有し、整列された単語群のうちの特定部分のみについて、単語もしくは単語とその重要度の組を抽出し、これを再利用可能な形で提示するようにしたことを特徴とする関連キーワード自動抽出装置。

【請求項2】 前記構成において、特定された部分集合Aに対して、これに含まれる部分集合Bが文書集合選定部により特定された場合に、部分集合Aに含まれる文書群に出現する単語の統計情報と、部分集合Bに含まれる文書群に出現する単語の統計情報との差分を、部分集合Bにおける各単語の重要度に加味することで、部分集合Bに出現する各単語の重要度を算出して単語ランキングに反映することを特徴とする請求項1に記載の関連キーワード自動抽出装置。

【請求項3】 文書集合選定部に各文書の重みを付与する機能を設け、特定された文書集合の各文書に含まれる単語の重要度に当該文書の重みを加味することにより当該単語の重要度を算出して単語ランキングに反映することを特徴とする請求項1または2に記載の関連キーワード自動抽出装置。

【請求項4】 対象文書集合全体において出現度合いが高頻度または低頻度である単語をあらかじめ定められた閾値を考慮して除外することにより、再利用の際に有効性の高い単語のみが選別できることを特徴とする請求項1乃至3のいずれかに記載の関連キーワード自動抽出装置。

【請求項5】 単語の長さなどその単語の特徴量に応じて除外のための閾値を変化させることにより再利用の際に有効性の高い単語のみが選別できることを特徴とする請求項4に記載の関連キーワード自動抽出装置。

【請求項6】 単語の出現位置や出現する文脈の情報を管理する出現情報管理部を有し、単語の重要度にその単語の出現情報の種類に応じてあらかじめ定められた重みを加味することにより当該単語の重要度を算出して単語ランキングに反映することを特徴とする請求項1乃至5のいずれかに記載の関連キーワード自動抽出装置。

【請求項7】 単語の品詞など、各単語の属性情報を管理する言語属性管理部を有し、当該単語の属性に応じてあらかじめ定められた重みを加味することにより当該単語の重要度を算出して単語ランキングに反映することを

特徴とする請求項1乃至6のいずれかに記載の関連キーワード自動抽出装置。

【請求項8】 抽出された単語同士、またはあらかじめ指定された単語群と抽出された単語との間の文字列としての包含関係を、定められた条件により判定する文字列包含関係判定部を有し、当該単語同士に文字列としての包含関係があると判定された場合に、指定された条件に従って、長単位の文字列のみ、もしくは短単位の文字列のみ、もしくは重要度の高い方の文字列のみ、もしくは短単位の文字列および長単位の文字列と短単位の文字列との差分の双方、のいずれかを選択することにより、再利用の際に有効性の高い単語のみが選別できることを特徴とする請求項1乃至7のいずれかに記載の関連キーワード自動抽出装置。

【請求項9】 単語の品詞など、各単語の属性情報を管理する言語属性管理部を有し、当該単語の属性や、指定された部分集合または文書全体における出現頻度、分布等を考慮することにより、抽出された単語を分類して提示できることを特徴とする請求項1乃至8のいずれかに記載の関連キーワード自動抽出装置。

【請求項10】 分類された単語群のそれぞれについて、その集合を代表する単語を付与する代表語付与部を設け、分類された単語群を代表する代表語群のみ、もしくは代表語と全ての単語を提示できることを特徴とする請求項9に記載の関連キーワード自動抽出装置。

【請求項11】 辞書を用いて対象文書集合の各文書に出現する単語または単語の組の出現頻度や分布などの統計情報があらかじめ抽出されている文書集合に対して、文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書の検索を行なう文書検索部と、文書検索部45において検索された文書について、入力された検索式と文書との間の適合度を計算する文書ランキング部46とを有して成り、文書ランキング部におけるランキング結果を関連キーワード自動抽出装置へ送付し、また関連キーワード自動抽出装置からフィードバックされた関連キーワードを検索条件入力部へ入力することが可能な文書検索装置。

【請求項12】 文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書の検索を行なう文書検索部とを有して成り、前記検索条件入力部は、ユーザが検索条件を入力するのに加えて関連キーワード自動抽出装置から送られてきた関連キーワードを検索条件として入力することが可能な文書検索装置。

【請求項13】 辞書を用いて対象文書集合の各文書に出現する単語または単語の組の出現頻度や分布などの統計情報があらかじめ抽出されている文書集合に対して、文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書

の検索を行なう文書検索部と、文書検索部 45 において検索された文書について、入力された検索式と文書との間の適合度を計算する文書ランキング部 46 とを有して成る文書検索装置と、

前記文書検索装置に接続された関連キーワード自動抽出装置とから構成され、

前記文書検索装置の文書ランキング部から出力されたランキング結果を関連キーワード自動抽出装置へ送付し、また関連キーワード自動抽出装置から文書検索装置の検索条件入力部へ関連キーワードをフィードバックしてキーワード検索を行なうようにしたことを特徴とする文書検索システム。

【請求項 14】 文書検索装置と関連キーワード自動抽出装置との間には文書集合選定部が設けられ、文書検索装置の文書ランキング部から出力されたランキング結果は文書集合選定部に送付されて文書の特定が行なわれ、前記関連キーワード自動抽出装置 48 には、文書集合選定部 47 が特定した文書の部分集合が入力されることを特徴とする請求項 13 記載の文書検索システム。

【請求項 15】 関連キーワード自動抽出装置には、請求項 1 乃至 10 のいずれかに記載の関連キーワード自動抽出装置が用いられることを特徴とする請求項 13 または 14 記載の文書検索システム。

【請求項 16】 文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書の検索を行なう文書検索部とを有して成る文書検索装置と、

前記文書検索装置に接続された関連キーワード自動抽出装置とから構成され、

前記文書検索装置の検索条件入力部は、ユーザが検索条件を入力するのに加えて関連キーワード自動抽出装置から送られてきた関連キーワードを検索条件として入力してキーワード検索を行なうようにしたことを特徴とする文書検索システム。

【請求項 17】 関連キーワード自動抽出装置には、請求項 1 乃至 10 のいずれかに記載の関連キーワード自動抽出装置が用いられることを特徴とする請求項 16 記載の文書検索システム。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、特定の文書集合から、その文書集合を特徴づける語句をキーワードとして抽出するための関連キーワード自動抽出装置、および前記関連キーワード自動抽出装置を利用した文書検索装置に関する。

【0002】

【従来の技術】文書検索装置において、ユーザが必要とする文書を得るためには、適切な検索語を利用した検索式を入力する必要があるが、ユーザ自身が適切な検索語を想起し難い、という問題がある。そこで従来、ユーザ

が入力した検索語に対して、関連語辞書などを利用して検索語に関連する語を提示することにより、ユーザの再検索を助ける手法などが取られてきた。しかしながら、こうした手法はあらかじめ静的にさだめられた関連語辞書の性質に依存するため、検索対象となる文書の特性に即した関連語が得られない。また、得られた単語で検索した結果少なくとも 1 件以上の文書が得られることが保証されない、という欠点があった。

【0003】

【発明が解決しようとする課題】本発明は前記の課題を解決するもので、特定された文書集合における各単語の出現頻度・分布などの統計情報と、検索対象文書全体における単語の統計情報とを考慮して単語の重要度を算出し、これにもとづいて単語をその重要度によってランキングし、ランクの一部である単語群を抽出することにより、実際の検索対象文書の特性に即し、かつ品質の高い関連キーワード群を高速かつ動的に抽出できる、関連キーワード自動抽出装置を提供することを目的とする。

【0004】また、前記関連キーワード自動抽出装置から得られた関連キーワード群を利用して検索を実行した場合、少なくとも 1 件以上の検索結果が得られることを保証する文書検索装置及びこれらを用いた文書検索システムを提供することを目的とするものである。

【0005】

【課題を解決するための手段】本発明は、上記目的を達成するため、関連キーワード自動抽出装置として、各文書に付与された属性情報やユーザが入力した検索式などに基づいて文書の部分集合を特定する文書集合選定部と、各単語の対象文書全体における統計情報や各文書ごとに出現する単語とその統計情報を管理する単語統計情報管理部と、各単語の全文書または各文書内統計情報を基に、特定された文書の部分集合に出現する各単語の重要度を算出して重要度の順に整列する単語ランキング部とを設け、単語統計情報管理部により、文書全体、および特定された文書部分集合における各単語の統計情報を高速に求めることが可能であり、特定された文書集合に出現する各単語を、その重要度の順に高速にランキングし、その一部を関連キーワードとして提示することができる。

【0006】さらに、前記構成に加えて、単語の属性情報や文書内の出現位置を管理する手段などを設けることにより、単語の重みを変化させ、あるいはランキング後の単語群から特定の条件を満たす単語を削除することで、抽出される単語群の関連語としての精度を向上させることができ、また、抽出された単語群を、語の属性や統計的性質により分類することで、よりわかりやすい関連キーワード提示を行なうことができる。

【0007】また本発明は、上記目的を達成するため、関連キーワード自動抽出装置と連携した文書検索装置を含む文書検索システムを構成し、抽出された関連キー

10

20

30

40

50

ードを入力として再利用することにより、抽出された関連キーワードが対象文書の特性に合ったものであり、かつ検索対象が同一の文書群であるならば、そのキーワードによって検索結果が少なくとも1件以上得られることが保障されるため、効率的かつ容易に再検索を行なうことができる。

【0008】

【発明の実施の形態】本発明の請求項1に記載の発明は、辞書を用いて対象文書集合の各文書に出現する単語または単語の組の出現頻度や分布などの統計情報があらかじめ抽出されている文書集合に対して、各文書に付与された属性情報やユーザが入力した検索式などに基づいて文書の部分集合を特定する文書集合選定部と、各単語の対象文書全体における統計情報、および各文書ごとの当該文書に出現する単語とその統計情報を管理する単語統計情報管理部と、各単語の全文書および各文書ごとの統計情報を基に、特定された部分集合に出現する各単語の重要度を算出して重要度の順に整列する単語ランキング部とを備えたものであり、整列された単語群のうちの特定部分のみについて、単語もしくは単語とその重要度の組を抽出し、これを再利用可能な形で高速提示するという作用を有する。

【0009】本発明の請求項2に記載の発明は、請求項1に記載の関連キーワード自動抽出装置において、特定された部分集合Aに対して、これに含まれる部分集合Bが文書集合選定部により特定された場合に、部分集合Aに含まれる文書群に出現する単語の統計情報と、部分集合Bに含まれる文書群に出現する単語の統計情報との差分を、部分集合Bにおける各単語の重要度に加味することで、部分集合Bに出現する各単語の重要度を算出して単語ランキングに反映するようにしたものである。

【0010】本発明の請求項3に記載の発明は、請求項1または2に記載の関連キーワード自動抽出装置において、文書集合選定部に各文書の重みを付与する機能を設け、特定された文書集合の各文書に含まれる単語の重要度に当該文書の重みを加味することにより当該単語の重要度を算出して単語ランキングに反映するようにしたものである。

【0011】本発明の請求項4に記載の発明は、請求項1乃至3のいずれかに記載の関連キーワード自動抽出装置において、対象文書集合全体において出現度合いが高頻度または低頻度である単語をあらかじめ定められた閾値を考慮して関連キーワード抽出の対象から除外することにより、再利用の際に有効性の高い単語のみが選別できるようにしたものである。

【0012】本発明の請求項5に記載の発明は、請求項4に記載の関連キーワード自動抽出装置において、単語の長さなどその単語の特徴量に応じて除外のための閾値を変化させることにより再利用の際に有効性の高い単語のみが選別できるようにしたものである。

【0013】本発明の請求項6に記載の発明は、請求項1乃至5のいずれかに記載の関連キーワード自動抽出装置において、単語の出現位置や出現する文脈の情報を管理する出現情報管理部を有し、単語の重要度にその単語の出現情報の種類に応じてあらかじめ定められた重みを加味することにより当該単語の重要度を算出して単語ランキングに反映するようにしたものである。

【0014】本発明の請求項7に記載の発明は、請求項1乃至6のいずれかに記載の関連キーワード自動抽出装置において、単語の品詞など、各単語の属性情報を管理する言語属性管理部を有し、当該単語の属性に応じてあらかじめ定められた重みを加味することにより当該単語の重要度を算出して単語ランキングに反映するようにしたものである。

【0015】本発明の請求項8に記載の発明は、請求項1乃至7のいずれかに記載の関連キーワード自動抽出装置において、抽出された単語同士、またはあらかじめ指定された単語群と抽出された単語との間の文字列としての包含関係を、定められた条件により判定する文字列包含関係判定部を有し、当該単語同士に文字列としての包含関係があると判定された場合に、指定された条件に従って、長単位の文字列のみ、もしくは短単位の文字列のみ、もしくは重要度の高い方の文字列のみ、もしくは短単位の文字列および長単位の文字列と短単位の文字列との差分の双方、のいずれかを選択することにより、再利用の際に有効性の高い単語のみが選別できるようにしたものである。

【0016】本発明の請求項9に記載の発明は、請求項1乃至8のいずれかに記載の関連キーワード自動抽出装置において、単語の品詞など、各単語の属性情報を管理する言語属性管理部を有し、当該単語の属性や、指定された部分集合または文書全体における出現頻度、分布等を考慮することにより、抽出された単語を分類して提示できるようにしたものである。

【0017】本発明の請求項10に記載の発明は、請求項9に記載の関連キーワード自動抽出装置において、分類された単語群のそれぞれについて、その集合を代表する単語を付与する代表語付与部を設け、分類された単語群を代表する代表語群のみ、もしくは代表語と全ての単語を提示できるようにしたものである。

【0018】本発明の請求項11に記載の発明は、文書検索装置として、辞書を用いて対象文書集合の各文書に出現する単語または単語の組の出現頻度や分布などの統計情報があらかじめ抽出されている文書集合に対して、文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書の検索を行なう文書検索部と、文書検索部において検索された文書について、入力された検索式と文書との間の適合度を計算する文書ランキング部とを備えたものであり、文書ランキング部におけるランキング結果を関連キ

ワード自動抽出装置へ送付し、また関連キーワード自動抽出装置からフィードバックされた関連キーワードを検索条件入力部へ入力するという作用を有する。

【0019】本発明の請求項 1 2 に記載の発明は、文書検索装置として、文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書の検索を行なう文書検索部とを備えたものであり、前記検索条件入力部は、ユーザが検索条件を入力するのに加えて関連キーワード自動抽出装置から送られてきた関連キーワードを検索条件として入力する

という作用を有する。

【0020】本発明の請求項 1 3 に記載の発明は、文書検索システムとして、辞書を用いて対象文書集合の各文書に出現する単語または単語の組の出現頻度や分布などの統計情報があらかじめ抽出されている文書集合に対して、文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書の検索を行なう文書検索部と、文書検索部において検索された文書について、入力された検索式と文書との間の適合度を計算する文書ランキング部とを有して成る文書検索装置と、前記文書検索装置に接続された関連キーワード自動抽出装置とを備えたものであり、前記文書検索装置の文書ランキング部から出力されたランキング結果を関連キーワード自動抽出装置へ送付し、また関連キーワード自動抽出装置から文書検索装置の検索条件入力部へ関連キーワードをフィードバックしてキーワード検索を行なうという作用を有する。

【0021】本発明の請求項 1 4 に記載の発明は、請求項 1 3 記載の文書検索システムにおいて、文書検索装置と関連キーワード自動抽出装置との間には文書集合選定部が設けられ、文書検索装置の文書ランキング部から出力されたランキング結果は文書集合選定部に送付されて文書の特定が行なわれ、前記関連キーワード自動抽出装置 4 8 には、文書集合選定部 4 7 が特定した文書の部分集合が入力されるようにしたものである。

【0022】本発明の請求項 1 5 に記載の発明は、請求項 1 3 または 1 4 記載の文書検索システムにおいて、関連キーワード自動抽出装置には、請求項 1 乃至 1 0 のいずれかに記載の関連キーワード自動抽出装置が用いられるようにしたものである。

【0023】本発明の請求項 1 6 に記載の発明は、文書検索に必要な条件式を入力する検索条件入力部と、入力された検索条件にしたがって対象文書集合から文書の検索を行なう文書検索部とを有して成る文書検索装置と、前記文書検索装置に接続された関連キーワード自動抽出装置とを備えたものであり、前記文書検索装置の検索条件入力部は、ユーザが検索条件を入力するのに加えて関連キーワード自動抽出装置から送られてきた関連キーワードを検索条件として入力してキーワード検索を行なうという作用を有する。

【0024】本発明の請求項 1 7 に記載の発明は、請求項 1 6 記載の文書検索システムにおいて、関連キーワード自動抽出装置には、請求項 1 乃至 1 0 のいずれかに記載の関連キーワード自動抽出装置が用いられるようにしたものである。

【0025】以下に、本発明の具体的な実施の形態について、添付の図面を参照して説明する。

【0026】(実施の形態 1) 最初に、本発明の第 1 の実施の形態について説明する。図 1 は本発明の第 1 の実施の形態に係る関連キーワード自動抽出装置の構成を示したブロック図である。まず、対象となる文書集合 1 1 に対し、辞書 1 2 を利用して、前処理として動作する統計情報抽出部 1 3 により、文書集合全体における単語の頻度・分布などの単語統計情報 1 4、および各文書ごとの当該文書に含まれる単語の統計情報である文書内単語統計情報 1 5 を抽出しておく。図 2 (a) は単語統計情報の構造を示すテーブル構成図であり、図 2 (b) は文書内単語統計情報の構造を示すテーブル構成図である。単語統計情報 1 4 は、統計情報抽出部 1 3 によって抽出された単語の統計情報を例えば図 2 (a) に示すようなテーブルとして格納する。このテーブルを利用することにより、例えば単語「インターネット」の全文書中での総出現頻度や出現文書数を高速に求めることができる。また、文書内単語統計情報 1 5 は各文書ごとの単語の統計情報を例えば図 2 (b) に示すようなテーブルとして格納する。これにより、例えば文書番号 0 0 1 0 には単語「インターネット」が 5 回、単語「WWW」が 2 回出現する、といった各文書ごとの統計情報を高速に求めることができる。

【0027】関連キーワード自動抽出装置 1 6 は、文書全体の単語統計情報 1 4 および文書内単語統計情報 1 5 を管理する単語統計情報管理部 1 7 と、単語の重要度を算出する単語ランキング部 1 8 と、対象文書の部分集合を特定する文書集合選定部 1 9 と、文書集合選定部 1 9 への選定条件を入力する手段である条件入力部 2 0 とから構成される。

【0028】かかる構成を有する関連キーワード自動抽出装置 1 6 の動作について以下説明する。最初に、条件入力部 2 0 に対して入力された条件により、文書集合選定部 1 9 が文書集合を特定する。文書集合は、次の 3 種類の手段のいずれかまたはその組み合わせにより特定される。

(1) 文書の属性により文書集合を特定する。この場合、文書集合選定部 1 9 は文書の所属するジャンルなど、文書にあらかじめ付与された属性値により文書を選択する手段を有し、条件入力部 2 0 により指定された属性値に合致する文書群を部分集合として採用する。

(2) 検索式により文書集合を特定する。この場合、文書集合選定部 1 9 は条件入力部 2 0 で入力された検索式に適合する文書を特定する文書検索手段を有し、これを

利用して検索の結果得られる文書群を部分集合として採用する。なおその際、文書検索手段に検索式との適合度を判定して文書を適合度の順にランキングする機能があるならば、検索結果のうちの特定部分、例えば上位 10 文書を部分集合として採用しても良い。

(3) ユーザにより指定された文書集合。この場合、文書集合選定部 19 は条件入力部 20 においてユーザが直接指定した (複数の) 文書を部分集合として採用する。

【0029】文書集合選定部 19 は、以上により選定された文書集合を各文書を一意に決定する識別子の集合、例えば文書番号のリストとして単語統計情報管理部 17 *

*に渡す。単語統計情報管理部 17 は、特定された文書集合に対して、文書ごとに文書番号から文書内単語統計情報 14 を調べ、当該文書に出現する単語とそれぞれの文書内の出現頻度を得る。次に得られた単語すべてについて単語統計情報 15 を調べ、当該単語の全文書における頻度や分布情報を得る。

【0030】ここで得られた各種統計情報は単語ランキング部 18 に渡され、各単語の重要度が算出される。ある単語 W の重要度 S (W) は、例えば次のようにして算出することができる。

【数 1】

$$S(W) = C * \sum_{j=0}^n \{TF_j(W) * IDF(W)\} * FN(W)$$

ただし

C : 定数

n : 特定された文書集合に含まれる文書数

TF_j(W) : 文書 D_j における単語 W の出現頻度

FN(W) : 特定された文書集合中で単語 W を含む文書数

である。

【0031】また IDF(W) は、単語 W の idf 値と呼ばれる指標であり、例えば以下の式により計算される。

$$IDF(W) = 1 - \log(DF(W) / N)$$

ただし、

DF(W) : 文書全体において単語 W が出現する文書数

N : 全文書数

である。

【0032】IDF(W) は、単語 W がより多くの文書に出現する (すなわちより一般的な語である) 場合にその値が小さくなる。これにより、対象文書全体において比較的良好に出現する語の重要度を低く抑えることができる。さらに FN(W) を考慮することで、特定された文書集合に多く現れる単語の重要度を高くでき、結果その特定文書集合に特徴的な語に高い重要度を与えることができる。なお、上記算出法において、TF(W) をその単語が含まれる文書の文書サイズ (文字数や含まれる単語の異なり数など) や単語の総出現頻度などで正規化してもよい。

【0033】単語ランキング部 18 は、特定された部分集合中の全文書に含まれる全単語について重要度計算を*

※行い、その後全単語を重要度の順に整列する。最後に、整列された単語群から特定部分、例えば上位 10 単語を採用し、単語、もしくは単語とその重要度の組として提示する。なお、抽出の際に重要度だけでなく、重要度算出に利用した各種統計情報などを同時に提示してもよい。また、抽出された関連キーワードとその重要度の組を、例えばユーザの履歴として蓄積していくこともできる。このようにすることにより、ユーザの興味の範囲や嗜好などをキーワードとその重みのベクトルとして表現することが可能となり、このベクトルを他の操作、例えば文書集合の検索に利用するなど、広い応用が可能である。

【0034】以上の計算式を利用すると、例えば図 3 に示す例のようにして関連キーワード自動抽出を行うことができる。この図 3 は関連キーワード自動抽出動作の処理手順の流れを示す図である。図 3 において、文書番号リスト 31 が入力された単語統計情報管理部 17 は、該当する文書番号 (例えば 0010、0341 等) に出現する単語およびその頻度を文書ごとに出力し、文書内単語統計情報 33、34、35 を得る。同時に、ここで求められたすべての単語に対して、全文書中での統計情報 32 を得る。次にこれらの統計情報 32、33、34、35 が単語ランキング部 18 に渡される。単語ランキング部 18 では、各種統計情報 32~35 を基に、例えば前記の式を利用して各単語の重要度を計算する。図 3 の場合だと、以下ようになる (ただし、C を 1、N を 10000 とする)。

$$IDF(\text{アプレット}) = 1 - \log(86 / 10000)$$

$$= 5.756$$

$$S(\text{アプレット}) = 2 * 5.756 + 6 * 5.756 * 2$$

$$= 92.096$$

$$IDF(\text{インターネット}) = 1 - \log(1129 / 10000)$$

$$= 3.181$$

$$S(\text{インターネット}) = (3 * 3.181 + 1 * 3.181 + 2 * 3$$

11

12

181)*3

$$\begin{aligned}
 &= 57.258 \\
 \text{IDF (CGI)} &= 1 - \log(79/10000) \\
 &= 5.840 \\
 \text{S (CGI)} &= (4 * 5.756) * 1 \\
 &= 23.024 \\
 \text{IDF (WWW)} &= 1 - \log(615/10000) \\
 &= 3.789 \\
 \text{S (WWW)} &= (5 * 3.789) * 1 \\
 &= 18.945 \\
 \text{IDF (JAVA)} &= 1 - \log(161/10000) \\
 &= 5.129 \quad 6 \\
 \text{S (JAVA)} &= (6 * 5.129 + 3 * 5.129 + 3 * 5.129) * 3 \\
 &= 184.644 \\
 \text{IDF (SUN)} &= 1 - \log(35/10000) \\
 &= 6.655 \\
 \text{S (SUN)} &= (6 * 6.655) * 1 \\
 &= 39.930 \\
 \text{IDF (スクリプト)} &= 1 - \log(813/10000) \\
 &= 3.510 \\
 \text{S (スクリプト)} &= (5 * 3.510) * 1 \\
 &= 17.550
 \end{aligned}$$

【0035】単語ランキング部18では以上のように求められた重要度により単語を整理し、整理後の単語リスト37を得る。ここで、ランキングされた単語の上位3語を抽出するという指定になっているとすれば、単語リスト37における上位3語である「JAVA」「アプレット」「インターネット」が関連キーワードとして抽出される。

【0036】以上では辞書に登録された一単語を抽出の対象としてきたが、一般に単語だけでなく、単語の組でもよい。単語の組とは、名詞の連続により構成される複合語や、助詞「の」で結ばれる名詞の組、助詞「を」

「が」で結ばれる名詞と動詞の組などを指す。これらの統計情報が単語と同様に事前に抽出できているのであれば、上記で示した手法がそのまま適用でき、単語の組を関連キーワードとして抽出することができる。

【0037】なお、関連キーワード入力装置16は、文書集合選定部19および条件入力部20を別構成としてもよい。特に文書集合選定部19が検索式による文書検索手段を有する場合には、後出の図7に示すような別構成とすることで、文書検索装置による文書番号を入力として受け、出力される関連キーワードを文書検索装置の検索式入力部に反映させることができる。

【0038】このように、本実施の形態によれば、対象となる文書のうちの一部である文書の部分集合が特定された際、当該部分集合に含まれる各文書に出現する各単語それぞれについて重要度を計算して重要度の順に整理し、整理された単語群のうちの一部を抽出して関連キー

ワードとすることで、動的かつ高速に対象となる文書の特性に即した関連キーワードを求めることができるという効果を持つ。

【0039】また、上記のようにして得られた関連キーワードは、同一文書を対象とする文書検索装置への入力として利用することができ、その場合、対象文書の特性にあった的確なキーワードを再利用できるだけでなく、当該関連キーワードは必ず対象文書に含まれることが保証されるため、これを利用して検索した場合に必ず検索結果が得られるという効果も持つ。

【0040】また、得られた関連キーワードを同一の対象文書集合または別の対象文書集合を対象とする文書検索装置への入力として利用することができ、その場合には、関連キーワード抽出の対象となった文書集合において特徴的であるキーワードをもとに、同一または別の文書集合を検索することができ、特に別の文書集合を検索対象とする文書検索装置の場合に、当該キーワードを異なった特性を持つ文書集合に対しても適用することができるという効果をもつ。

【0041】また、抽出されたキーワードをユーザに提示して選択させるという構成とすることで、ユーザが再検索を実行する際、キーボードから再度検索条件を入力する代わりに、関連キーワードを、例えばマウスのクリックなど単純な操作で選択することが可能となり、再検索における操作を軽減して検索の効率を高めると同時に、検索の操作に不慣れなユーザでも簡単に利用できるという効果を持つ。

30

40

50

【0042】また、抽出された関連キーワードにその重要度も付加して提示することにより、例えば検索条件との適合度を計算して文書をランキングするような文書検索装置において、検索条件中の各単語に対して重みを付与することができる文書検索装置であれば、抽出されたキーワードとその重要度をそのまま入力とすることで、より高精度の検索結果を得ることができるという効果を持つ。

【0043】また、抽出された関連キーワードとその重要度の組を、例えばユーザの履歴として蓄積していくことにより、ユーザの興味の範囲や嗜好などをキーワードとその重みのベクトルとして表現することが可能となり、このベクトルを他の文書集合の検索に利用するなど、広い応用が可能であるという効果も持つ。

【0044】（実施の形態2）次に、本発明の第2の実施の形態について実施の形態1に示したブロック図と同じ図1を利用して説明する。この第2の実施の形態では、文書集合選定部19が2種類の文書集合Aおよび文書集合Bを特定する。ここで、文書集合Bは文書集合Aの部分集合となっている。例えば、ある検索式で検索を行った結果得られる文書集合Aと、そのうちで関連する文書群としてユーザが指定した文書集合Bとが特定される場合や、文書の属性により特定された文書集合Aと、その中でさらに検索式により絞り込まれた文書集合Bとが特定される場合などである。

【0045】この場合、例えば以下の式により算出される単語の分布指標を当該単語の重要度に乗算するなどの手法により、単語の重要度を算出する。

$$DI(A, B, W) = \{(NA/DA(W)) * (D * S2(Aプレット) = 92.096 * \{(100/10) * (2/3)\} = 613.973$$

$$S2(Aプレット) = 92.096 * \{(100/10) * (2/3)\} = 613.973$$

$$S2(インターネット) = 57.258 * \{(100/28) * (3/3)\} = 204.493$$

$$S2(CGI) = 23.024 * \{(100/9) * (1/3)\} = 85.274$$

$$S2(WWW) = 18.945 * \{(100/14) * (1/3)\} = 45.107$$

$$S2(JAVA) = 184.644 * \{(100/20) * (3/3)\} = 923.220$$

$$S2(SUN) = 39.930 * \{(100/5) * (1/3)\} = 266.200$$

$$S2(スクリプト) = 17.550 * \{(100/10) * (1/3)\} = 58.500$$

$$= 58.500$$

* B(W) / NB) }

ただし、

DA(W) : 部分集合Aにおける単語Wの出現する文書数

DB(W) : 部分集合Bにおける単語Wの出現する文書数

NA : 部分集合Aの総文書数

NB : 部分集合Bの総文書数

【0046】これは、部分集合Bにおいて高い頻度で出現し、かつ部分集合Aにおける出現頻度が低いもののほど高い値となる。上式において高い値となる語は部分集合Aにおいて部分集合Bの弁別性に大きく寄与するものであり、部分集合Bをより特徴づけるキーワードであるといえる。例えば、図3に示す例において、文書番号リスト31が部分集合Bであるとし、これを含む部分集合A（総文書数100とする）も同時に指定されている場合で、部分集合A中の各単語の出現文書数が以下の通りであるとすると、

$$DA(Aプレット) = 10$$

$$DA(インターネット) = 28$$

$$DA(CGI) = 9$$

$$DA(WWW) = 14$$

$$DA(JAVA) = 20$$

$$DA(SUN) = 5$$

$$DA(スクリプト) = 10$$

【0047】この場合各単語の重要度S2(W)は、実施の形態1で説明した各単語の重要度S(W)に各単語の重みDI(A, B, W)を乗算した値となり、以下のように計算される。

となり、重要度の順に整列すると

S 2 (J A V A)	= 9 2 3. 2 2 0
S 2 (アプレット)	= 6 1 3. 9 7 3
S 2 (S U N)	= 2 6 6. 2 0 0
S 2 (インターネット)	= 2 0 4. 4 9 3
S 2 (C G I)	= 8 5. 2 7 4
S 2 (スクリプト)	= 5 8. 5 0 0
S 2 (WWW)	= 4 5. 1 0 7

の順となる。したがって、上位 3 語を関連キーワードとして抽出するのであれば、「J A V A」「アプレット」「S U N」が関連キーワードとなる。

【0048】上記の計算式は一例であり、部分集合 B において高い頻度で出現し、かつ部分集合 A における出現頻度が低いものほど高い値となるような他の計算式を利用してもよい。

【0049】このように、本実施の形態によれば、特定された 2 種類の部分集合間における頻度分布の差異を考慮することにより、より高精度な関連キーワードを得ることができるという効果を持つ。

【0050】（実施の形態 3）次に、本発明の第 3 の実施の形態について実施の形態 1 に示したブロック図と同じ図 1 を利用して説明する。この第 3 の実施の形態では、文書集合選定部 19 に各文書の重みを付与する機能を設ける。例えば、ユーザが文書を指定する場合に、各文書に対して関連度を指標として 5 段階の評価値を与える場合や、検索式による検索の結果得られる文書が検索式との適合度によりランキングされている場合に 1 位に 10 点、2 位に 9 点、といった重みを与える場合などである。単語ランキング部は各文書に付与された重みを、当該文書に含まれる単語に対して、例えば乗算するなどして加味し重要度算出を行う。なお、各文書に与える重みは負の値であってもよい。例えば、ユーザが文書を特定する際、関連する文書には 2 点、まったく関連しない文書には -1 点を与える、という重み付与も許す。これにより、関連する文書にも関連しない文書にも含まれる（かつあまり一般的でない）語の重要度を低くすることができる。

【0051】このように、本実施の形態によれば、特定した文書集合に含まれる各文書に対して重みを与えることにより、より重要な文書に含まれる単語ほど高い重要度となるような計算式とすることで、文書それぞれの重要度を勘案した高精度な関連キーワードが得られるという効果を持つ。

【0052】（実施の形態 4）次に、本発明の第 4 の実施の形態について説明する。図 4 は本発明の第 4 の実施の形態に係る関連キーワード自動抽出装置のブロック図である。この第 4 の実施の形態では、第 1 の実施の形態の構成に加えて閾値設定部 22 を有して成り、この閾値設定部は単語統計情報管理部 17 との間でデータの送受ができるようになっている。また、この実施の形態にお

いては、単語統計情報管理部 17 には閾値による単語除外機能が付与されている。かかる構成において、単語統計情報管理部 17 は各単語の統計情報を出力する際、あらかじめ定められた閾値設定 22 を参照し、極端に高頻度または低頻度の単語はその場で候補から除外して単語ランキング部 18 に当該単語の情報を出力しない構成とすることができる。例えば、閾値 1 を「全文書の 50% 以上に出現する単語」と設定し、閾値 2 を「1 文書にしか出現しない単語」と設定することで、これらの単語が重要度算出に与える悪影響を事前に防ぐことができ、かつ処理の高速化を図ることができる。

【0053】なおその際、単語の長さなど当該単語の特徴量に応じて、閾値を何種類かに設定してもよい。例えば、日本語の場合で「二文字以上の語は全体の 50% 以上、一文字の語は全体の 30% 以上」といった閾値設定を行うことで、各語の特性にあわせて除外する単語の範囲を設定する。

【0054】このように、本実施の形態によれば、対象文書集合全体において出現度合いが高頻度または低頻度である単語をあらかじめ定められた閾値を考慮して除外することにより、キーワード抽出処理を高速化でき、かつ再利用の際に有効性の高い単語のみが選別できるという効果を持つ。

【0055】（実施の形態 5）次に、本発明の第 5 の実施の形態について説明する。図 5 は本発明の第 5 の実施の形態に係る関連キーワード自動抽出装置の構成を示すブロック図である。この第 5 の実施の形態に係る関連キーワード自動抽出装置は、第 1 の実施の形態において説明したような、文書全体の単語統計情報 14 および文書内単語統計情報 15 を管理する単語統計情報管理部 17、単語ランキング部 18、対象文書の部分集合を特定する文書集合選定部 19、および文書集合選定部 19 への選定条件入力手段である条件入力部 20 を有する基本構成に加えて、単語ランキング部 18 と連動して単語の属性などの各種情報を利用することにより、抽出される関連キーワード群の質を向上させることを目的とするものである。図 5 において、符号 25 は出現情報管理部、26 は単語属性情報管理部、27 は文字列包含関係判定部であり、これらの機能部は関連キーワード自動抽出装置 29 に含まれて単語ランキング部と連動する。また 28 は代表語付与部であり、この代表語付与部 28 は単語ランキング部 18 からデータを受けて関連キーワードを

出力する。また、関連キーワード自動抽出装置 2 9 に対して、外部機能部として、対象文書集合 1 1 からのデータを基に単語が出現する位置の情報を抽出する単語出現位置情報抽出部 2 3 が設けられ、この単語出現位置情報抽出部 2 3 からは出現位置情報 2 4 が出力される。この出現情報は出現情報管理部 2 5 へ送付される。

【0 0 5 6】かかる構成を有する本発明の第 5 の実施の形態について、その動作を説明する。この実施の形態の動作においては、まず対象となる文書集合 1 1 に対し、辞書 1 2 を利用して、前処理として動作する統計情報抽出部 1 3 により、対象文書集合 1 1 全体における単語の出現頻度・分布などの単語統計情報 1 4、および各文書ごとの当該文書に含まれる単語の統計情報である文書内単語統計情報 1 5 を抽出しておく。同時に、必要があれば単語位置情報抽出部 2 3 により、単語の出現位置情報 2 4 も抽出しておく。図 6 は単語出現位置情報抽出部 2 3 によって抽出された出現位置情報 2 4 のデータ構造の一例を表すテーブル構成図である。出現位置情報は例えば図 6 に示すようなテーブルとして格納される。各文書ごとにその文書に出現する単語と出現位置（例えば文書の先頭からのバイトオフセット）、出現区分などが格納される。

【0 0 5 7】そして関連キーワード自動抽出動作に際しては、各単語に対して出現情報管理部 2 5 に問い合わせを行い、当該単語の出現位置や出現文脈などの情報を得、これを重要度算出に加味する。例えば、検索対象とする文書すべてが、タイトル（または見出し）、サブタイトル、本文、といった要素から構成されている文書である場合、当該単語がこれら要素のうちいずれに含まれているかによって、

タイトルに含まれる場合には 3 点

サブタイトルに含まれる場合には 2 点

本文に含まれる場合には 1 点

といったような「重み」を各単語の重要度に乗算する、といった手法で重要度を算出する。

【0 0 5 8】あるいは、出現位置の情報を利用してもよい。例えば部分集合が検索式により特定される場合で、この検索式に含まれる単語が参照可能である場合、検索式に含まれる単語と、現在重要度計算の対象となっている単語との間の文字数が、

2 文字以内なら 3 点

1 0 文字以内なら 2 点

1 0 文字以上なら 1 点

といったような「重み」を当該単語の重要度に乗算する、といった手法で重要度を算出することも可能である。

【0 0 5 9】また、本実施の形態の別の態様として、各単語に対して、単語属性情報管理部 2 6 に問い合わせを行い、当該単語の品詞や分類など、その単語の属性を得、これを重要度算出に加味する。例えば、当該単語の

品詞に着目し、

固有名詞ならば 5 点

普通名詞ならば 4 点

形容詞、形容動詞ならば 2 点

動詞、副詞ならば 1 点

その他自立語でないもの（助詞、助動詞など）ならば 0 点

といったような「重み」を各単語の重要度に乗算する、といった手法で重要度を算出することも可能である。

【0 0 6 0】また、本実施の形態の別の態様として、ある 2 つの単語間の文字列としての包含関係を判定する文字列包含関係判定部 2 7 を用いて、抽出された単語同士、もしくはあらかじめ指定された単語群のうちの単語と抽出された単語との間に包含関係があるか否かを判定し、包含関係があると判定された場合に、抽出する単語を制限する。ここであらかじめ指定された単語群とは、例えば部分集合の特定に検索式を利用した場合の検索式に含まれる単語などである。包含関係の判定においては、あらかじめ定められた設定により、以下の判定基準のいずれか一つ（または一つ以上）を満たす場合を包含関係と認定することができる。

（1）単語 A と単語 B とが前方において一致しかつ単語 A が単語 B より短い場合、（2）単語 A と単語 B とが後方において一致しかつ単語 A が単語 B より短い場合、

（3）単語 A が単語 B の部分でありかつ前方、後方ともに一致しない場合、（4）単語 A と単語 B との関係が

（1）～（3）のいずれかを満たし、かつ単語 B の構成要素と完全に一致する場合、

【0 0 6 1】例えば、（1）の基準では「東京都」に対する「東京」が部分語と判定される。以下、同様にし、（2）の基準では「新発売」に対する「発売」が、（3）の基準では「大感謝祭」に対する「感謝」が、それぞれ部分語と判定される。（4）の基準は、英語における部分語判定の際に重要であり、この基準に従えば "artificial intelligence" に対して "art" や "tell" は部分語とはならないが、"artificail" や "intelligence" は部分語と判定される。

【0 0 6 2】上記基準により、部分語関係にあると判定された 2 つの語について、そのどちらを関連キーワードとして採用するかについても、以下のいずれかの基準（あらかじめ設定されているものとする）に従う。

（1）長単位の単語を採用する

（2）短単位の単語を採用する

（3）重要度の高い単語を採用する

（4）短単位の単語および長単位の単語と短単位の単語との差分を採用する

【0 0 6 3】例えば、単語「東京都」が重要度 1 0 で、単語「東京」が重要度 7 でそれぞれ抽出され、かつ両者に部分語関係が成立した場合、（1）の基準に従うと文字列として長い「東京都」が採用され、（2）の基準に

従うと文字列として短い「東京」が採用され、(3)の基準に従うとより重要度の高い「東京都」が採用されることになる。(4)の基準は、例えば単語 "artificial intelligence" と "artificial" との間に部分語関係が成立した場合に、"artificial" および "intelligence" を関連キーワードとして採用するものであり、主に英語文書において効果的である。

【0064】あらかじめ指定された単語群との間に部分語関係が成立する単語の場合、(3)以外の手法が利用できる。その場合、「短単位（もしくは長単位）であれば関連キーワードとして採用しない」といった処理となる。抽出された単語同士に部分語関係が成立する場合には、いずれの手法も利用可能である。

【0065】また、本実施の形態の別の態様として、抽出された関連キーワード群を、各語の属性や統計情報を利用して分類して提示する。語の属性として品詞を利用すると、例えば固有名詞とそれ以外に分類して提示することができる。あるいは、語の属性としてシソーラス辞書を利用し、各語をシソーラスにおける分類に対応する形で分類して提示することも可能である。また、統計情報を利用した分類とは、例えば特定された文書集合における各語の出現文書数により分類する手法などがあげられる。その場合、例えば「出現文書数が文書集合の8割以上であるか否か」といった基準で分類することで、その語が再検索に利用される際の絞り込みの効果を事前に確認することができる。なお、分類にあたり語の属性としてシソーラス辞書を利用する場合、分類された単語群に対して、シソーラスの上位ノードに相当する語を代表語として与え、単語群をその語で代表させることも可能である。同様に、単語の統計情報14を利用する場合には、分類された単語群において、例えば最も出現頻度の高い語を代表語として採用してもよい。

【0066】このように、本実施の形態によれば、単語が出現した位置の情報を利用することで、文書の構造や単語間の距離の情報を考慮した関連キーワードの抽出が行なえ、高精度な関連キーワード抽出が可能となるという効果を持つ。

【0067】また、単語の品詞など、各単語の属性情報を考慮することにより、各属性の特徴に応じた関連キーワードの抽出が行なえ、高精度な関連キーワード抽出が可能となるという効果を持つ。

【0068】また、単語間の文字列としての包含関係を考慮することにより、同じような意味や用途である単語を排除して関連キーワードの抽出が行なえ、関連キーワード全体としての冗長性を抑えることができるという効果を持つ。

【0069】また、抽出された関連キーワードを分類し、必要があれば各分類に対応する代表語を設定することで、抽出されたキーワードの一覧性や傾向、再利用の際の有効性などをあらかじめ確認して関連キーワードの

抽出が行なえ、関連キーワードとしての使いやすさを向上することができるという効果を持つ。

【0070】（実施の形態6）次に、本発明の第6の実施の形態について説明する。図7は本発明の第6の実施の形態に係る文書検索装置の構成およびこれと関連キーワード自動抽出装置とを組み合わせる実現した文書検索システムの構成を示すブロック図である。この文書検索装置41は、前記第1、第2、第3、第4または第5の実施の形態に係る関連キーワード自動抽出装置と連携して動作するものである。

【0071】本実施形態における文書検索装置41は、文書検索に必要な条件式を入力する検索条件入力部44と、入力された検索条件にしたがって文書の検索を行なう文書検索部45と、文書検索部45において検索された文書について入力された検索式と文書との間の適合度を計算する文書ランキング部46とを有して成る。この文書検索装置41は、連携して動作する関連キーワード自動抽出装置48と同一の対象文書集合11を検索対象とするものであり、単語統計情報抽出に利用するのと同じ辞書12を利用して、あらかじめ索引生成部42により作成された文書検索用の索引43を利用して検索を行う。また、本実施形態における関連キーワード自動抽出装置48は、文書集合選定部47を別構成としたものであり、関連キーワード自動抽出装置48には、文書集合選定部47が特定した文書の部分集合の各要素に対応する文書の識別子の集合（一意である文書番号のリストなど）が入力される。

【0072】以上の構成を備えた本実施の形態について、その動作を説明する。最初に検索条件入力部44に入力された検索条件をもとに、文書検索部45が検索用索引43を参照して検索条件に適合する文書を特定する。ここで得られた文書集合をそのまま検索結果文書50としてもよいが、さらに文書ランキング部46において、入力された検索式と文書との間の適合度を計算して適合度の高い順に文書を整列したものを検索結果とする、といった構成にしてもよい。こうして得られた検索結果の文書集合50は、ユーザに検索結果として返すのと同時に、文書集合選定部47に渡される。文書選定部47では、文書ランキング部46から渡された文書集合のすべてまたは一部を関連キーワード自動抽出装置48への入力として採用する。文書が適合度の順にランキングされているのであれば、検索結果の文書集合のうち例えば上位10文書を選定する、という構成にしてもよい。また、あらかじめ文書ごとに付与された属性情報を利用できるのであれば、これを利用して特定の属性値を持つ文書のみを選定する、という構成にしてもよい。

【0073】文書集合選定部47により特定された文書の部分集合は関連キーワード自動抽出装置48に送られ、前記第1、第2、第3、第4または第5の実施の形態に示したような手順で関連キーワード群49を抽出す

る。こうして得られた関連キーワード群49は検索条件入力部44に戻され、ユーザに提示される。ユーザは提示された関連キーワード群から必要なものを選択して新たな検索条件とし、再度検索を実行させることができる。これにより、本実施の形態によれば、関連キーワード自動抽出装置によって上記のようにして得られた関連キーワードは、同一文書を対象とする文書検索装置への入力として利用することができ、その場合、対象文書の特性にあった的確なキーワードを再利用できるだけでなく、当該関連キーワードは必ず対象文書に含まれることが保証されるため、これを利用して検索した場合に必ず検索結果が得られるという効果も持つ。

【0074】（実施の形態7）次に、本発明の第7の実施の形態について説明する。図8は本発明の第7の実施の形態に係る文書検索装置の構成およびこれと関連キーワード自動抽出装置とを組み合わせることで実現した文書検索システムの構成を示すブロック図である。この文書検索装置51は、第6の実施の形態に係る文書検索装置41と同様、前記第1、第2、第3、第4または第5の実施の形態に係る関連キーワード自動抽出装置と連携して動作するものである。

【0075】本実施形態における文書検索装置51は、文書検索に必要な条件式を入力する検索条件入力部54と、入力された検索条件にしたがって文書の検索を行なう文書検索部55とを有して成る。本実施形態における文書検索装置51は、連携して動作する関連キーワード自動抽出装置52とは異なる対象文書集合56を検索対象とするものであり、文書検索部55が対象文書集合56に接続される構成となっている。なお、その検索手法についての詳細は問わない。

【0076】以上の構成を備えた本実施形態における動作について、以下説明する。最初に指定された条件にしたがって関連キーワード自動抽出装置52が動作し、関連キーワード群53を出力する。文書検索装置51における検索条件入力部54は、関連キーワード群53を入力としてユーザに提示し、ユーザは提示された関連キーワードのうち必要なもののみを選択して、検索対象となる対象文書集合56に対する検索を実行し、検索結果文書57を得ることができる。

【0077】このように、本実施の形態によれば、関連キーワード自動抽出装置52によって得られた関連キーワードを同一の対象文書集合または別の対象文書集合を対象とする文書検索装置51への入力として利用することができ、その場合には、関連キーワード抽出の対象となった文書集合において特徴的であるキーワードをもとに、同一または別の文書集合を検索することができ、特に別の文書集合を検索対象とする文書検索装置の場合に、当該キーワードを異なった特性を持つ文書集合に対しても適用することができるという効果をもつ。

【0078】

【発明の効果】以上説明したように、本発明によれば、関連キーワード自動抽出装置を、文書の部分集合を特定する文書集合選定部と、対象文書全体または個々の文書ごとに出現する単語とその統計情報を管理する単語統計情報管理部と、文書の部分集合に出現する各単語の重要度を算出して重要度の順に整列する単語ランキング部とにより構成したため、文書全体、および特定された文書部分集合における各単語の統計情報を高速に求めることが可能であり、特定された文書集合に出現する各単語を、その重要度に基づいて高速にランキングし、その一部を関連キーワードとして提示することができる。

【0079】また、前記構成に加えて、単語の属性情報や文書内の出現位置を管理する手段などを設けることにより、単語の重みを変化させ、あるいはランキング後の単語群から特定の条件を満たす単語を削除することで、抽出される単語群の関連語としての精度を向上させることができる。また、抽出された単語群を、語の属性や統計的性質により分類することで、よりわかりやすい関連キーワード提示を行なうことができる。

【0080】さらに、関連キーワード自動抽出装置と連携した文書検索装置を含む文書検索システムを構成し、抽出された関連キーワードを入力として再利用することにより、抽出された関連キーワードが対象文書の特性に合ったものであり、かつ検索対象が同一の文書群であるならば、そのキーワードによって検索結果が少なくとも1件以上得られることが保障されるため、効率的かつ容易に再検索を行なうことができる等の効果が得られる。

【図面の簡単な説明】

【図1】本発明の第1乃至第3の実施の形態に係る関連キーワード自動抽出装置の構成を示すブロック図

【図2】（a）前記実施の形態における単語統計情報の構造を示すテーブル構成図

（b）前記実施の形態における文書内単語統計情報の構造を示すテーブル構成図

【図3】前記実施の形態における関連キーワード自動抽出動作の処理手順の流れを示す図

【図4】本発明の第4の実施の形態に係る関連キーワード自動抽出装置の構成を示すブロック図

【図5】本発明の第5の実施の形態に係る関連キーワード自動抽出装置の構成を示すブロック図

【図6】前記実施の形態における単語出現位置情報抽出部によって抽出された出現位置情報のデータ構造の一例を表すテーブル構成図

【図7】本発明の第6の実施の形態に係る文書検索装置の構成構成およびこれと関連キーワード自動抽出装置とを組み合わせることで実現した文書検索システムの構成を示すブロック図

【図8】本発明の第7の実施の形態に係る文書検索装置の構成構成およびこれと関連キーワード自動抽出装置とを組み合わせることで実現した文書検索システムの構成を示すブ

ック図

【符号の説明】

11、56 対象文書集合

12 辞書

13 統計情報抽出部

14 単語統計情報

15 文書内単語統計情報

16、29、48、52 関連キーワード自動抽出装置

17 単語統計情報管理部

18 単語ランキング部

19 文書集合選定部

20 条件入力部

21、49、53 関連キーワード群

22 閾値設定

23 単語出現位置情報抽出部

* 24 出現位置情報

25 出現情報管理部

26 単語属性情報管理部

27 文字列包含関係判定部

28 代表語付与部

41、51 文書検索装置

42 索引生成部

43 検索用索引

44、54 検索条件入力部

10 45、55 文書検索部

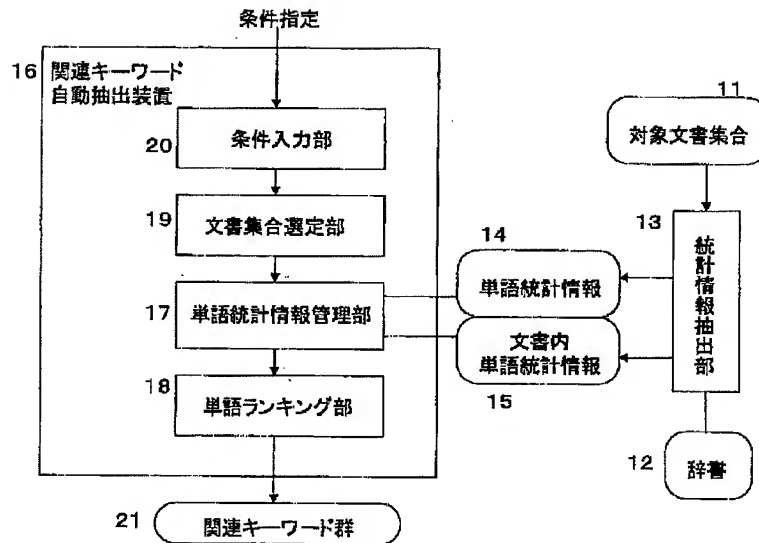
46 文書ランキング部

47 文書集合選定部

50、57 検索結果文書

*

【図1】



【図2】

14

(a)

単語統計情報

単語	総出現 頻度	出現 文書数

インターネット	1026	542
インターハイ	15	10
インターバル	2078	1129
インタビュー	104	91
インタフェース	5288	2275

15

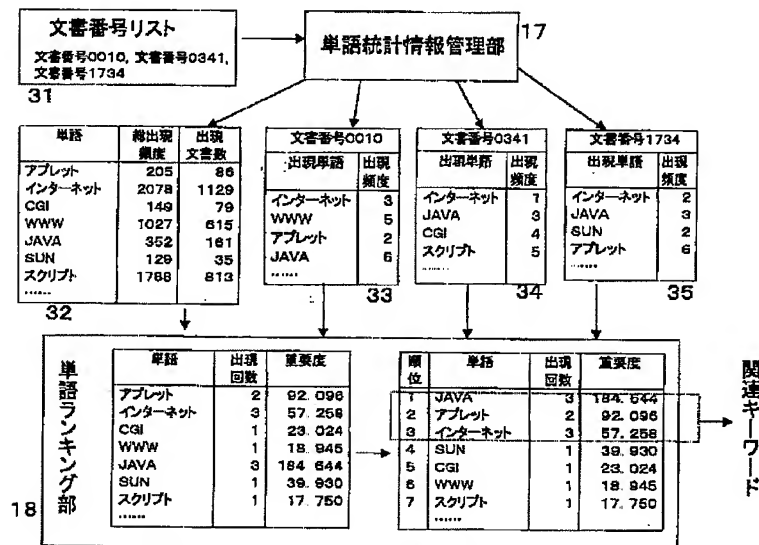
(b)

文書内単語統計情報

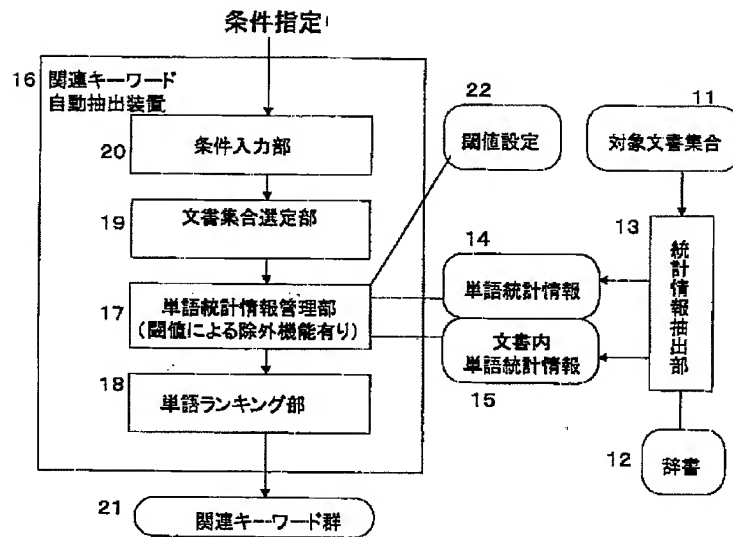
文書番号	出現単語	出現 頻度

文書0010	インターネット	5
文書0010	WWW	2
文書0011	イントラネット	6
文書0011	LAN	10
文書0011	WAN	2

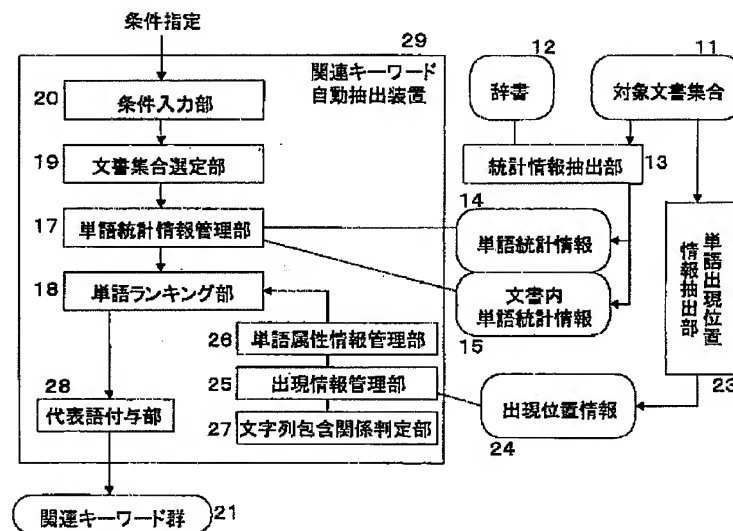
【図3】



【図 4】



【図 5】



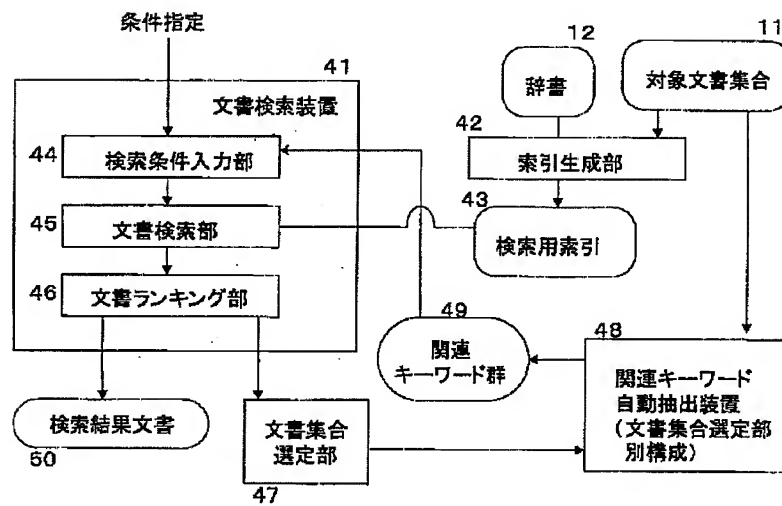
【図6】

24

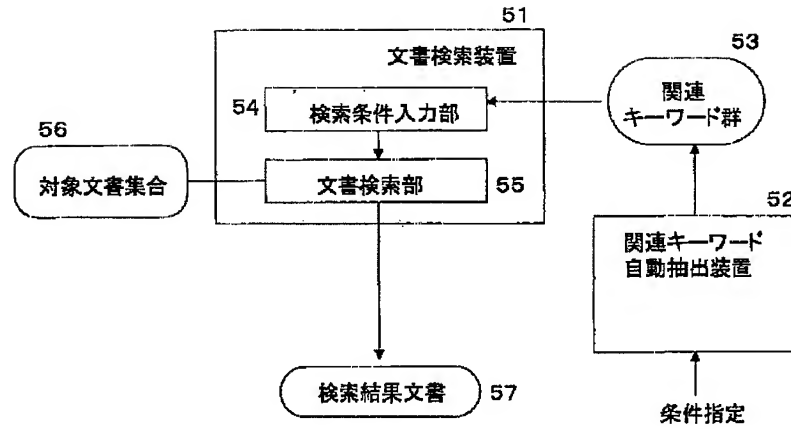
単語出現位置情報

文書番号	出現単語	出現位置	出現区分
.....			
文書0010	インターネット	10368	タイトル
文書0010	WWW	10384	タイトル
文書0010	近年	10390	本文
文書0010	インターネット	10396	本文
文書0010	成長	10412	本文
.....			

【図7】



【図8】



フロントページの続き

(72)発明者 野 本 昌 子
 大阪府門真市大字門真1006番地 松下電器
 産業株式会社内

(72)発明者 稲 葉 光 昭
 大阪府門真市大字門真1006番地 松下電器
 産業株式会社内
 (72)発明者 福 重 貴 雄
 大阪府門真市大字門真1006番地 松下電器
 産業株式会社内